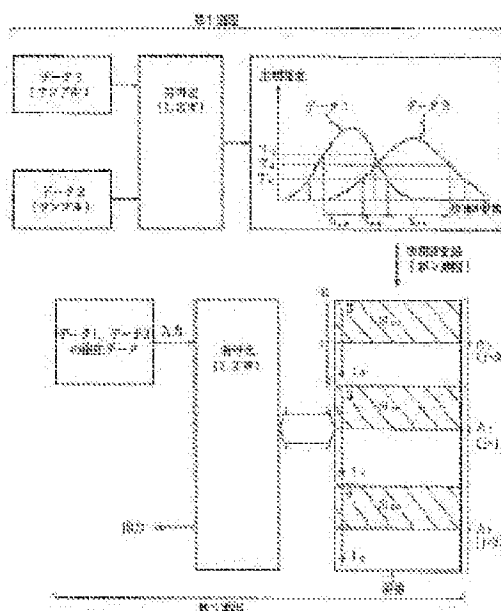


METHOD FOR COMPRESSING PLURAL KINDS OF DATA

Patent number: JP6028149 (A)
Publication date: 1994-02-04
Inventor(s): YOSHIDA SHIGERU; OKADA YOSHIYUKI; NAKANO YASUHIKO; CHIBA HIROTAKA +
Applicant(s): FUJITSU LTD +
Classification:
- **international:** G06F5/00; G06T9/00; H03M7/30; H03M7/46; G06F5/00; G06T9/00; H03M7/30; H03M7/46; (IPC1-7): G06F15/66; G06F5/00; H03M7/30
- **european:**
Application number: JP19920183288 19920710
Priority number(s): JP19920183288 19920710

Abstract of JP 6028149 (A)

PURPOSE:To provide a method to obtain high compressibility without increasing the processing time of the data compressing for plural kinds of data which compresses and encodes input data by using dynamic dictionary type algorithm of LZW codes by uniting the dictionary retrieval of encoding wherein character string which are investigated as to plural kinds of data and frequently appear are initially registered. **CONSTITUTION:**Partial strings whose appearance frequency detected by the LZW encoding of sample data 1 and 2 become high in common to plural kinds of data are extracted as a common partial string group S00 and initially registered in a dictionary area A0. Further, partial strings which have high frequencies by the data are extracted as characteristic partial string groups S10 and S20 and initially registered in dictionary areas A1 and A2. When an input character string wherein the data 1 and 2 are mixed is encoded, a registered partial string which matches the input character string to the longest length is retrieved in a dictionary and the input character string is encoded by using the group number (j) of a partial string group belonging to the retrieved partial string and the registration number (i) of the retrieved character string in the partial string group.



Data supplied from the **espacenet** database — Worldwide

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平6-28149

(43) 公開日 平成6年(1994)2月4日

(51) Int.Cl. ⁵	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 5/00	H	9189-5B		
15/66	3 3 0 H	8420-5L		
H 0 3 M 7/30		8522-5J		

審査請求 未請求 請求項の数6 (全 23 頁)

(21) 出願番号 特願平4-183288

(22) 出願日 平成4年(1992)7月10日

(71) 出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中1015番地

(72) 発明者 吉田 茂

神奈川県川崎市中原区上小田中1015番地

富士通株式会社内

(72) 発明者 岡田 佳之

神奈川県川崎市中原区上小田中1015番地

富士通株式会社内

(72) 発明者 中野 泰彦

神奈川県川崎市中原区上小田中1015番地

富士通株式会社内

(74) 代理人 弁理士 竹内 進 (外1名)

最終頁に続く

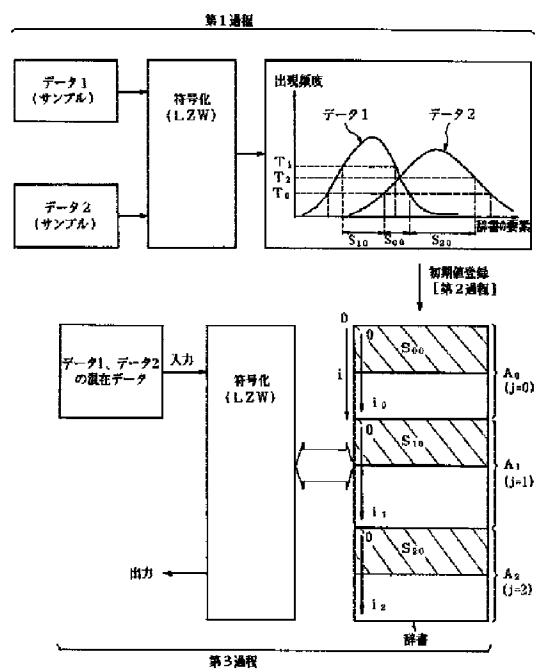
(54) 【発明の名称】 複数種類データのデータ圧縮方法

(57) 【要約】

【目的】 L Z W符号等の動的辞書型アルゴリズムを用いて入力データを圧縮符号化する複数種類データのデータ圧縮方法に関し、複数種類のデータを対象に調べた高頻度に出現する文字列を初期登録した符号化における辞書探索を一本化し、処理時間を増加させることなく高い圧縮率が得られるようにする。

【構成】 サンプルデータ1、2のL Z W符号化で検出した出現頻度が、複数種類のデータで共通に高頻度となる部分列を共通部分列群 S_{00} として抽出して辞書領域 A_0 に初期登録する。また複数種類のデータ毎に高頻度となる部分列を固有部分列群 S_{10} 、 S_{20} として抽出して辞書領域 A_1 、 A_2 に初期登録する。データ1、2が混在する入力文字列を符号化する際には、入力文字列に最長一致する登録済み部分列を辞書から検索し、検索した部分列の属する部分列群の群番号 j と部分列群内の検索文字列の登録番号 i とを用いて入力文字列を符号化する。

本発明の原理説明図



1

2

【特許請求の範囲】

【請求項1】複数種類のデータが混在する入力データを符号化して圧縮する複数種類データのデータ圧縮方法に於いて、

複数種類の各データを相異なる部分列に分けて辞書に登録し、各データ毎に入力文字列に最長一致する辞書に登録済みの部分列を検索し、検索した部分列の登録番号で表わして入力文字列を圧縮符号化し、該符号化における部分列の出現頻度を検出する第1過程と、

前記第1過程で検出した出現頻度が、複数種類のデータで共通に高頻度となる部分列を共通部分列群 (S_{00}) として抽出し、該共通部分列群 (S_{00}) に特定の群番号 ($j=0$) を付けて辞書領域 (A_0) を確保し、該辞書領域 (A_0) に該共通部分列群 (S_{00}) に属する各部分列を初期登録し、また前記第1過程で検出した出現頻度が、複数種類のデータ毎に高頻度となる部分列を固有部分列群 (S_{10} , S_{20}) として抽出し、各固有部分列群 (S_{10} , S_{20}) 毎に特定の群番号 ($j=1, 2$) を付けて辞書領域 (A_1 , A_2) を確保し、該辞書領域 (A_1 , A_2) に該当する固有部分列群 (S_{10} , S_{20}) に属する各部分列を初期登録する第2過程と、

複数種類のデータが混在する入力文字列を符号化する際に、入力文字列に最長一致する登録済み部分列を前記辞書から検索し、検索した部分列の属する部分列群 (S_{00} , S_{10} , S_{20}) の群番号 ($j=0, 1, 2$) と該部分列群内での検索文字列の登録番号 (i_1) とを用いて入力文字列を符号化する第3過程と、を備えたことを特徴とする複数種類データのデータ圧縮方法。

【請求項2】複数種類のデータが混在する入力データを符号化して圧縮する複数種類データのデータ圧縮方法に於いて、

複数種類の各データを相異なる部分列に分けて辞書に登録し、各データ毎に入力文字列に最長一致する辞書に登録済みの部分列を検索し、検索した部分列の登録番号で表わして入力文字列を圧縮符号化し、該符号化における部分列の出現頻度を検出する第1過程と、

前記第1過程で検出した出現頻度が複数種類のデータで共通に高頻度となる共通部分列群 (S_{00}) と各データ毎に高頻度となる部分列群とを合わせた固有部分列群 (S_{10} , S_{20}) を抽出し、各固有部分列群 (S_{10} , S_{20}) 毎に特定の群番号 ($j=1, 2$) を付けて辞書領域 (A_1 , A_2) を確保し、該辞書領域 (A_1 , A_2) に前記共通部分列群 (S_{00}) を合せた各固有部分列群 (S_{10} , S_{20}) に属する各部分列を初期登録する第2過程と、

複数種類のデータが混在する入力文字列を符号化する際に、入力文字列に最長一致する部分列を前記辞書から検索し、検索した部分列が前記共通部分列群 (S_{00}) に属するときは該共通部分列群 (S_{00}) 内の登録番号 (i) を用いて符号化し、一方、固有部分列群 (S_{10} , S_{20}) に属するときは、該固有部分列群の群番号 (A_1 , A_2) と該群内の登録番号 (i_1) とを用いて符号化する第3過程と、を備えたことを特徴とする複数種類データのデータ圧縮方法。

【請求項3】請求項1, 2記載の複数種類データのデータ圧縮方法に於いて、前記第3過程では、符号化すべき入力文字列に最長一致する部分列が、前回符号化で最長一致した部分列と同じ部分列群に属しているときは、該群内の登録番号のみを用いて符号化し、前回符号化で最長一致した部分列と異なる部分列群に属しているときは、該群番号と該群内の登録番号を用いて符号化することを特徴とする複数種類データのデータ圧縮方法。

【請求項4】請求項1, 2記載の複数種類データのデータ圧縮方法に於いて、前記第1過程では、出現する全種類のデータのサンプルごとに符号化を行って相異なる部分列の出現頻度を計数することを特徴とする複数種類データのデータ圧縮方法。

【請求項5】請求項1, 2記載の複数種類データのデータ圧縮方法に於いて、前記第3過程にあっては、入力文字列を辞書の最長一致する部分列の検索で符号化した際に、該符号化済み文字に次の入力文字を加えた文字列を、符号化文字列が属する部分列群に新たな参照番号を付けて登録することを特徴とする複数種類データのデータ圧縮方法。

【請求項6】請求項1, 2記載の複数種類データのデータ圧縮方法に於いて、前記第2過程にあっては、各部分列群ごとに最大登録個数を予め定めて該部分列群を登録するメモリ領域を割り当てておき、前各群に属する部分列の登録番号を各メモリ領域の先頭からの位置で表すことを特徴とする複数種類データのデータ圧縮方法。

【請求項7】請求項1, 2記載の複数種類データのデータ圧縮方法に於いて、前記第2過程にあっては、各部分列群ごとに最大登録個数を予め定めて該部分列群を登録するメモリ領域を割り当てておき、前各群に属する部分列の登録番号を各メモリ領域の先頭からの位置で表すことを特徴とする複数種類データのデータ圧縮方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は、LZW符号等の動的辞書型アルゴリズムを用いて入力データを圧縮符号化する複数種類データのデータ圧縮方法に関する。近年、文字コード、ベクトル情報、画像など様々な種類のデータがコンピュータで扱われるようになっており、扱われるデータ量も急速に増加してきている。

【0002】大量のデータを扱うときは、データの中の冗長な部分を省いてデータ量を圧縮することで、記憶容量を減らしたり、速く伝送したりできるようになる。様々なデータを1つの方式でデータ圧縮できる方法としてユニバーサル符号化が提案されている。ここで、本発明の分野は、文字コードの圧縮に限らず、様々なデータに適用できるが、以下では、情報理論で用いられている呼称を踏襲し、データの1ワード単位を文字と呼び、データが任意ワードつながったものを文字列と呼ぶことにする。

【0003】ユニバーサル符号の代表的な方法として、ジブレンベル (Ziv-Lempel) 符号がある (詳しくは、例えば、宗像『Ziv-Lempelのデータ圧縮法』、情報処

理、Vol. 26, No. 1, 1985年を参照のこと）。

ジブーレンベル符号では

①ユニバーサル型（スライド辞書型）と、

②増分分解型（Incremental parsing；動的辞書型）

の2つのアルゴリズムが提案されている。

【0004】さらに、ユニバーサル型アルゴリズムの改良として、LZSS符号、(T.C. Bell, "Better OPM/L Text Compression", IEEE Trans. on Commun., Vol. CO M-34, No. 12 Dec. 1986 参照)。また、増分分解型アルゴリズムの改良としては、LZW (Lempel-Ziv-Welch) 符

号がある (T.A. Welch, "A Technique for High-Performance Data Compression", Computer, June 1984 参照)。

【0005】これらの符号の内、高速処理ができること

と、アルゴリズムの簡単さからLZW符号が記憶装置の

ファイル圧縮などで使われるようになっていく。

【0006】従来の技術】従来のLZW符号の符号化アルゴリズム

を図15のフローチャートに示す。LZW符号化は、書き

替え可能な辞書をもち、入力文字列を相異なる部分列

に分け、この部分列を出現した順に参照番号を付けて辞

書に登録するとともに、現在入力している文字列を辞書

に登録してある最長一致する部分文字列の参照番号だけ

で表して、符号化するものである。

【0007】尚、増分分解型符号およびLZW符号の技

術は、特開昭59-231683号、米国特許第4, 5

58, 302号、米国特許第4, 814, 746号で開

示されている。図15の符号化処理は次のようになる。

ステップS1；予め全文字につき一文字からなる文字列

を初期値として登録してから符号化を始める。辞書の登

録数nを文字種数Aと置く。

【0008】カーソルをデータの先頭の位置に置く。

ステップS2；カーソルの位置からの文字列に一致する

辞書登録の最長文字列Sを見つける。

ステップS3；文字列Sの識別番号を $[10g_2 n]$ ビ

ットで表して出力する。但し、 $[x]$ はx以上の最小の

整数である。辞書登録数nを一つインクリメントする。

【0009】ステップS4；文字列Sにカーソルの最初

の文字Cを付加した文字列SCを辞書に登録する。カー

ソルはSの後の文字に移動させる。

ステップS5；入力データの終了をチェックし、終了す

るまでステップS2～S4の処理を繰り返す。

図16は従来のLZW復号化のアルゴリズムを示したフ

ローチャートであり、図15の符号化アルゴリズムと逆

の操作を行って入力符号から文字列を復元し、同時に辞

書を作成する。

【0010】このような従来のLZW符号では、複数の

異なる性質をもつ複数種類のデータが混在するデータ

を符号化すると、複数種のデータに合わせた辞書が作成

されて符号化が行なわれる。複数種が混在するデータとし

ては、例えば、文字コードと画像が混在するデータが挙げられる。辞書のサイズが十分大きいときは、出現した全てのデータ種を含む辞書が作られるため、個々のデータ種単独で圧縮した場合に比べて圧縮率が悪化するという問題がある。

【0011】また、辞書のサイズが小さいために1種類のデータ分しか登録できないようなときは、各データの種類の著しく異なれば、辞書をクリアして再学習するため、個々のデータ種に合わせた辞書が作られ、圧縮率は低下しない。しかし、データ中に同じデータ種が交互に出現するときは、そのつど学習し直すため、圧縮率が高められないという問題点があった。

【0012】この問題点を解決するため、本願発明者らは、データの種類ごとに辞書を分けて作成することで高い圧縮率を得るようにした方法を提案している。図17にデータの種類毎に辞書を作成して符号化するLZW符号化アルゴリズムを示す。図17のLZW符号化は次のようになる。

【0013】ステップS1；データの種類ごとに高頻度で出現する文字列を求め、辞書の初期値とする。

ステップS2；データの種類ごとに初期値を分割辞書に設定する。カーソルを1とし、辞書アドレス n_1 をデータjの初期値の個数 A_j とし、直前辞書番号を $pp=0$ とする。

【0014】ステップS3；カーソルをセットした位置からの入力文字列に一致する各辞書j中の最長の文字列 $S_j=S_1, S_2, \dots, S_k$ を見つける。

ステップS4；ステップS3で見つけた文字列 S_j の中から最長の文字列 S_p を求める。

【0015】ステップS5；現在辞書番号pと直前辞書番号ppが一致するかどうか判定する。

ステップS6；辞書番号が不一致の場合は、 $[10g_2 n_{pp}]$ ビットを用いて辞書が変わったことを示す識別番号0を表わし、また $[10g_2 K]$ ビットを用いて変わった辞書番号pを表わして出力する。

【0016】ステップS7；ステップS6の出力が済みまたはステップS5で辞書番号が一致した場合に、ステップS4で検索した文字部分列 S_p の番号を $[10g_2 n_p]$ ビットを用いて表わし、出力する。辞書アドレス n_p を1つインクリメントする。

ステップS8；文字列Sの次の文字をCにセットする。符号化済み文字列 S_p に文字Cを加えた文字列 S_pC を、辞書アドレス n_p で辞書に登録する。現在辞書番号pを直前辞書番号ppに置き替える。

【0017】カーソルを文字列Sの位置の文字に移動させる。

ステップS9；データ終了の有無を判別し、終了していなければステップS3に戻り、終了していれば一連の処理を終る。

図18は図17のLZW符号化アルゴリズムの変形を示

したもので、図16のステップS5、S6で行っている参照辞書が変化を示す情報の符号化出力を除いており、他の点は同じになる。

【0018】この図17、図18に示す複数種類データのデータ圧縮方法では、データの種類ごとに高頻度で出現する文字列を調べて、データの種類の辞書D1に予め設定しておき、複数の辞書D1から検索した最長一致文字列の中から最も一致長が長い文字列の辞書を選んで符号化するものである。このため高頻度の初期値を元にデータの種類が分類され、データ種に適する辞書D1が選ばれるため、高圧縮率を得ることができる。

【0019】

【発明が解決しようとする課題】しかしながら、複数種類のデータが混在するデータを図17、図18の方法で符号化する場合、高圧縮率は得られるものの、複数の辞書について最長一致する文字列を検索しなければならず、辞書検索に時間がかかるという問題がある。この辞書検索の問題は、ハードウェアで並列処理を行うようにすれば単一辞書を用いた従来のLZW符号化と同等の処理速度が得られるが、ソフトウェアによるシーケンシャル処理では辞書の複数の個数分の検索時間がかかり、処理速度が低下する問題があった。

【0020】本発明は、このような問題点を鑑みてなされたもので、複数種類のデータを対象に調べた高頻度に出現する文字列を初期登録した場合の符号化における辞書探索を一本化し、処理時間を増加させることなく高い圧縮率が得られるようにした複数種類データのデータ圧縮方法を提供することを目的とする。

【0021】

【課題を解決するための手段】図1は本発明の原理説明図である。まず本発明は、複数種類のデータが混在する入力データを符号化して圧縮する複数種類データのデータ圧縮方法として、次のようにする。

〔第1過程〕複数種類のデータをLZW符号化し、この符号化における辞書に登録した文字列の出現頻度を検出する。

【0022】〔第2過程〕第1過程のLZW符号化で検出した出現頻度が、複数種類のデータで共通に高頻度となる部分列を共通部分列群 S_{00} として抽出し、共通部分列群 S_{00} に特定の群番号 $j=0$ を付けて辞書領域 A_0 を確保し、この辞書領域 A_0 に共通部分列群 S_{00} に属する各部分列を初期登録する。

【0023】また第1過程のLZW符号化で検出した出現頻度が、複数種類のデータ毎に高頻度となる部分列を固有部分列群 S_{10} 、 S_{20} として抽出し、各固有部分列群 S_{10} 、 S_{20} 毎に特定の群番号 $j=1$ 、 $j=2$ を付けて辞書領域 A_1 、 A_2 を確保し、辞書領域 A_1 、 A_2 に該当する固有部分列群 S_{10} 、 S_{20} に属する各部分列を初期登録する。

【0024】〔第3過程〕複数種類のデータが混在する

入力文字列を符号化する際に、入力文字列に最長一致する登録済み部分列を辞書10から検索し、検索した部分列の属する部分列群の群番号 j と該部分列群内の検索文字列の登録番号 i とを用いて入力文字列を符号化する。

【0025】また本発明の他の複数種類データのデータ圧縮方法としては、辞書領域を各データに共通な領域を各データに固有な領域と一緒にしてもよい。この場合の処理は、次のようになる。

〔第1過程〕複数種類のデータをLZW符号化し、この符号化における辞書に登録した文字列の出現頻度を検出する。

【0026】〔第2過程〕第1過程のLZW符号化で検出した出現頻度が複数種類のデータで共通に高頻度となる共通部分列群 S_{00} と各データ毎に高頻度となる部分列群 S_{10} 、 S_{20} とを合わせた固有部分列群 $(S_{00}+S_{10}, S_{00}+S_{20})$ を抽出し、各固有部分列群 $(S_{00}+S_{10}, S_{00}+S_{20})$ 毎に特定の群番号 $j=1, 2$ を付けて辞書領域 A_1 、 A_2 を確保し、該辞書領域に各固有部分列群に属する各部分列を初期登録する。

【0027】〔第3過程〕複数種類のデータが混在する入力文字列を符号化する際に、入力文字列に最長一致する部分列を辞書から検索し、検索した部分列が共通部分列群 S_{00} に属するときは共通部分列群内の登録番号 n を用いて符号化し、一方、固有部分列群 S_{10} 、 S_{20} に属するときは、固有部分列群の群番号 $j=1, 2$ と該群内の登録番号 n_1 とを用いて符号化する。

【0028】

【作用】このような本発明の複数種類データのデータ圧縮方法にあっては、予め出現するデータ種が例えば2種類の場合を例にとると、辞書に初期登録する文字列を2種類のデータについて高頻度で共通に現れる文字列と、各種類のデータだけに高頻度で現れる文字列との3つの集合 S_{00} 、 S_{10} 、 S_{20} （文字列群）に分けて辞書に格納しておき、符号化時に入力データと最長一致する辞書中の文字列を、対応する集合ごとの参照番号 $j=0, 1, 2$ を付けて符号化する。

【0029】これにより複数種類のデータ毎に調べた高頻度の文字列を単一の辞書に初期登録して複数種類のデータが混在するデータの符号化を行うことができ、ソフトウェアによる符号化であっても辞書検索を通常のLZWと同等の処理速度で行って高い圧縮率を得ることができる。

【0030】

【実施例】図2は複数種類のデータが混在したデータを符号化する本発明の複数種類データのデータ圧縮方法を実現する装置構成の一実施例を示した実施例構成図である。図2において、16はCPUであり、CPU16に対してはプログラムメモリ18とデータメモリ30が接続される。プログラムメモリ18にはコントロールソフ

ト20、符号化ソフト22、初期値作成手段としての機能を備えた辞書作成ソフト14、出現頻度カウントテーブル26及び出現頻度格納テーブル28が設けられる。

【0031】符号化ソフト22は入力文字列に最長一致する辞書中の文字列を検索して辞書の参照番号を符号データとして出力する、例えばLZW符号化アルゴリズムを基本的に備える。また、復号化ソフト24は符号化ソフト22により符号化された入力符号列で辞書中の参照番号を検索し、対応する文字列を復号する例えばLZW復号化アルゴリズムを行う。

【0032】辞書作成ソフト14は符号化あるいは復号化に先立って行う初期値作成処理と符号化及び復号化の処理中に新たな文字列を辞書に登録する処理の2つを行う。この辞書作成ソフト14における初期値作成機能はデータメモリ30に格納された複数種類のデータを代表するサンプルデータ、例えば2種類のサンプルデータ1、2を対象に符号化ソフト22に従った符号化を行い、この符号化時に辞書から文字を検索して符号データとして出力する毎に、符号データとして検索された辞書中の文字列の参照番号の使用回数を出現頻度カウントテーブル26を使用してカウントアップし、文字列の出現頻度を検出する。

【0033】サンプルデータの符号化が終了したならば出現頻度カウントテーブル26の中のサンプルデータ1、2毎の出現頻度を参照し、2種類のデータ1、2に共通する高頻度の文字列の集合を初期値 S_{00} として登録し、またデータ1、2のそれぞれで独立に高頻度に生ずる文字列の集合を初期値 S_{10} 、 S_{20} として登録する。一方、データメモリ30には辞書10とデータバッファ32の各メモリ領域が確保される。

【0034】初期値作成時にはデータバッファ32には初期値作成の対象となる複数種類のサンプルデータ、例えばサンプルデータ1、2が格納され、また辞書10には初期値作成のための符号化時に辞書作成ソフト14で作成された文字列が参照番号と共に登録される。初期値作成が済むと、辞書10には辞書作成ソフト14で作成された複数種類、例えば2種類のデータ1、2の共通の初期値 S_{00} 、各データ1、2毎の初期値 S_{10} 、 S_{20} がそれぞれに割り当てられた領域 A_0 、 A_1 、 A_2 毎に初期登録が行われる。そしてデータバッファ32には新たに符号化しようとするデータ1、2が混在した文字列あるいは復号化しようとする符号列が格納され、符号化ソフト22による文字列の復号化あるいは復号化ソフト24による文字列の復元が行われる。

【0035】図3は本発明において2種類のサンプルデータ1、2を対象に辞書に登録する初期値の作成処理を示したフローチャートである。図3の初期値作成処理にあっては、まずステップS1で2種類のデータ1、2を対象にLZW符号化を行い、データ1、2に共通な高頻度の初期値 S_{00} を作成する。

【0036】続いてステップS2で共通の初期値 S_{00} を辞書の A_{00} 領域に格納してデータ1のみをLZW符号化し、データ1に特有の高頻度の初期値 S_{10} を作成する。続いてステップS3で共通の初期値 S_{00} を辞書の A_{00} 領域に格納してデータ2のみをLZW符号化し、データ2に特有の高頻度の初期値 S_{20} を作成する。具体的には、図4に示すようにサンプルデータ1、2を対象にLZW符号化を行って辞書に符号化済み文字列の部分列を参照番号と共に登録し、サンプルデータ1の符号化における出現頻度 f_1 とサンプルデータ2の符号化における出現頻度 f_2 のそれぞれを計数する。

【0037】図5はサンプルデータ1、2の符号化で得られた出現頻度を縦軸にとり辞書の要素番号(参照番号)を横軸にとって示した説明図である。図5において、サンプルデータ1、2中の要素(文字列)の出現頻度をそれぞれ f_1 、 f_2 とし、また共通初期値 S_{00} の閾値を T_0 、各サンプルデータ1、2特有の閾値を T_1 、 T_2 とすると、初期値 S_{00} 、 S_{10} 、 S_{20} の集合は次のようになる。

初期値 S_{00} ： $f_1 > T_0$ かつ $f_2 > T_0$ の要素の集合

初期値 S_{10} ： $f_1 \geq T_0$ かつ $f_2 \leq T_0$ かつ $f_1 > T_1$ の要素の集合

初期値 S_{20} ： $f_1 \leq T_0$ かつ $f_2 \geq T_0$ かつ $f_2 > T_2$ の要素の集合

このようにデータ1、2が混在した場合の符号化で作成される辞書要素の全体をデータ1、2に共通の集合 S_{00} とデータ1、2に固有の集合 S_{10} 、 S_{20} に分類して辞書に初期登録しておけば、この初期登録した辞書を用いた符号化で符号化中のデータが最長一致する辞書の参照番号がどの集合に属するかを調べることでデータ1、2の変移区間 S_{00} なのか特定データ1、2の区間 S_{10} または S_{20} にあるかを判別することができ、単一の辞書を用いてデータの種類に対応した効率の良い符号化を行うことができる。

【0038】図6は図3のステップS1に示したデータ1、2に共通の初期値 S_{00} を作成する初期値作成処理を詳細に示したフローチャートである。図6において、まずステップS1にデータ1、2のそれぞれにおける全ての単一文字を初期値として登録してから符号化を始める。また辞書の登録数 n を文字種数 A と置き、カーソルをデータの先頭位置にセットし、更に出現頻度を計数するカウンタ f_1 を変移要素 N 個分準備して0に初期化する。

【0039】次にステップS2でサンプルデータ1の入力を開始し、ステップS3でデータ入力終了をチェックした後、ステップS4に進んでカーソル位置からの文字列に一致する辞書中の最長の文字列 S を見付ける。続いてステップS5で見付けた最長一致の文字列に含まれる全てのセット文字列について出現頻度 f_1 を1つイン

クリメントする。

【0040】次にステップS6で辞書アドレスnを1つインクリメントし、符号化した最長一致文字列Sの次の文字をCとし、この次の1文字を文字列Sに加えた文字列SCを参照番号nを付けて辞書に登録する。そして、カーソルを文字列Sの次の文字に移動させ、ステップS2で次のサンプルデータ1を入力する。ステップS2～S6の処理の繰返しにより、ステップS3でサンプルデータ1の入力終了が判別されるとステップS7に進み、再びカーソルを1にセットし、サンプルデータにおける出現頻度計数のため、サンプルデータ2の全要素分N個のカウンタf₂を0にリセットし、ステップS8でサンプルデータ2の入力を開始する。

【0041】続いてステップS9を介してステップS10に進み、サンプルデータ2のカーソル位置からの文字列に一致する辞書中の最長一致する文字列Sを見付け、ステップS11で見付けた最長一致の文字列に含まれる全てのセット文字列について出現頻度f₂を1つインクリメントする。続いてステップS12で辞書番号nを1つインクリメントし、検索した最長一致文字列Sの次の1文字をCとし、最長一致文字列S₁に次の1文字Cを加えた文字列SCを参照番号nを付けて辞書に登録し、カーソルを文字列Sの次の文字に移動させ、再びステップS8に戻る。

【0042】ステップS9でサンプルデータ2の入力終了が判別されるとステップS13に進み、サンプルデータ1の出現頻度f₁及びサンプルデータ2の出現頻度f₂が共に閾値T₀となる辞書中の文字列を取り出して初期値S₀₀とする。図7は図3の初期値作成処理におけるステップS2及びS3の詳細を示したフローチャートである。

【0043】図7にあっては、まずステップS1で図6で作成したサンプルデータ1、2に共通な高頻度の初期値S₀₀を辞書Dに格納し、カーソルを1に合わせ、辞書アドレスをn₀及びn₁にセットし、出現頻度を計数するカウンタfを0にリセットする。続いてステップS2でまずサンプルデータ1を入力し、ステップS3を介してステップS4で辞書中の最長一致する文字列Sを見つけ、ステップS4で最長一致文字列に含まれる全てのセット文字列について出現頻度fを1つインクリメントする。

【0044】続いてステップS6で辞書アドレスnを1つインクリメントし、最長一致文字列Sの次の1文字をCとし、この1文字を最長一致文字列Sに加えた文字列SCに参照番号nを付けて辞書に登録する。続いてカーソルを文字列Sの後ろの1文字に移動させ、ステップS2に戻って次のサンプルデータ1を入力する。以上のステップS2～S6の処理の繰返しによりサンプルデータ1の符号化が済むとサンプルデータ2の符号化に切り替わり、同様な処理を繰返す。

【0045】ステップS3でデータ入力終了が判別されるとステップS7に進み、サンプルデータ1、2毎に計数されている出現頻度f₁、f₂に付き、閾値T₁、T₂以上となる辞書中の文字列を取り出してサンプルデータ1、2に特有な初期値S₁₀、S₂₀とする。図8はサンプルデータ1、2から作成された初期値S₀₀、S₁₀、S₂₀を用いた本発明によるLZW符号化アルゴリズムを示したフローチャートである。

【0046】図8において、まずステップS1において予めサンプルデータ1、2から作成した初期値S₀₀、S₁₀、S₂₀をそれぞれ辞書の領域A₀₀、A₁₀、A₂₀に格納する。また、各領域における既存の辞書登録の要素数n₀、n₁、n₂を設定する。図9は図8のLZW符号化で使用される辞書構成を示した説明図である。図9において、各符号は次の内容を示す。

A₀、A₁、A₂：共通部分、データ種1、データ種2の格納領域

N_{0max}、N_{1max}、N_{2max}：格納領域A₀、A₁、A₂の各格納領域の最大要素数

S₀₀、S₁₀、S₂₀：共通部分、データ種1、データ種2の初期値

A₀₀、A₁₀、A₂₀：共通部分、データ種1、データ種2の初期値の格納領域

n₀₀、n₁₀、n₂₀：共通部分、データ種1、データ種2の初期値の要素数

A₀₁、A₁₁、A₂₁：共通部分、データ種1、データ種2の既登録要素の格納領域

n₀、n₁、n₂：共通部分、データ種1、データ種2の既登録要素数

A₀₂、A₁₂、A₂₂：共通部分、データ種1、データ種2の空き領域

例えば、データ1、2に共通な高頻度をもつ初期値S₀₀を登録した辞書領域A₀についてみると、初期格納領域A₀₀に初期要素数n₀₀の初期値S₀₀を登録している。この領域A₀₀に続いて初期値S₀₀を用いた符号化で新たに登録された要素を含む既登録要素領域A₀₁が設けられ、ここまでの既存の登録要素数をn₀としている。また領域A₀、A₁、A₂については最大要素数をN_{0max}、N_{1max}及びN_{2max}と予め定めている。

【0047】再び図8を参照するに、辞書に対する初期登録が済むとステップS2で入力データと最長一致する辞書中の文字列(要素)を探索し、参照番号iを求める。続いて参照番号iが含まれる辞書領域A_iより辞書領域番号jを求める。この実施例ではデータは2種類ではあることから辞書領域jは図9に示すようにA₀、A₁、A₂の3つであり、辞書領域番号jはj=0、1、2のいずれかとなる。

【0048】次にステップS4で前回の辞書領域番号と今回求めた辞書領域番号jとが等しいか否かチェックし、等しければステップS6に進み、参照番号iを辞書

領域 A_1 に対応する番号 i_1 に変換して符号化出力する。一方、前回の辞書領域番号が今回求めた辞書領域番号 j に一致しなかった場合にはステップS5で新たな辞書領域番号 j を符号化してからステップS6で参照番号 i の符号化出力を行う。

【0049】続いてステップS7で辞書領域 A_1 に空きがあれば、その辞書領域の辞書アドレス n_1 を1つインクリメントし、最長一致した文字列に次の1文字を付加した文字列を辞書領域 A_1 に参照番号 n_1 を付けて追加登録する。ステップS8ではデータ終了の有無をチェックしており、データが終了しなければステップS2に戻って同様な処理を繰り返し、データが終了すれば一連の符号化処理を終わる。

【0050】図8のステップS6における参照番号 i を辞書領域 A_1 に対応する参照番号 i_1 に変換する処理は次のモード1～3に従って行う。

【モード1】

$0 \leq i < N_{0max}$; 辞書領域番号=0

辞書領域の対応番号 $i_0 = i + N_p$

【モード2】

$N_{0max} \leq i < N_{1max}$; 辞書領域番号=1

辞書領域の対応番号 $i_1 = i - N_{0max} + N_p$

【モード3】

$N_{1max} \leq i$; 辞書領域番号=2

辞書領域に対応番号 $i_2 = i - (N_{0max} + N_{1max}) + N_p$

ここで、 N_p は予約語の数であり、この実施例では例えば $N_p = 5$ の予約語を辞書領域の先頭に設けている。例えば、図10に示すように、辞書の先頭アドレス0～4を予約語領域とし、この辞書アドレス即ち参照番号0～4を各予約領域に示した意味をもつ情報として使用する。

【0051】即ち、参照番号0は辞書領域番号 A_0 を示し、参照番号1は辞書領域 A_1 を示し、また参照番号2は辞書領域番号 A_2 を示す。また、参照番号3は辞書の初期化を指令する制御コマンドとしての意味をもつ。更に参照番号4は符号化データの終了を示すEOF等を用いる。このため、実際の辞書領域は予約語領域に続くアドレス5、即ち参照番号5から開始され、モード1～3に示すように検索した参照番号 i に予約語数 N_p を加えることで実アドレスが求まる。

【0052】またモード1～3における辞書領域に対応した参照番号 i_0 、 i_1 、 i_2 は図9に示した辞書領域 A_0 、 A_1 、 A_2 における領域内での相対位置を示している。このため、絶対位置を示す参照番号 i に対し各領域 $A_0 \sim A_2$ 内での相対位置を示す参照番号 i_0 、 i_1 、 i_2 に変換することで、より少ない数値の参照番号とでき、符号化データのビット長を短縮して圧縮率を高めることができる。

【0053】また、前記モード1～3に示すようにして

求めた各領域の対応番号 $i_j = i_0, i_1, i_2$ は各領域の要素数 $n_j = n_0, n_1, n_2$ を用いて表現し得る最小ビット数である

$\lceil \log_2 (n_j + N_p) \rceil$ ビット

で符号化する。但し、 $\lceil X \rceil$ は X 以上の最小の整数を示している。

【0054】更に図8のステップS5における辞書領域番号 j の符号化にあつては、辞書領域番号 j を

$\lceil \log_2 (n_k + N_p) \rceil$ ビット

で符号化することになる。図11は図8のLZW符号化で得られた符号化データの説明図であり、図11にあつては符号化に使用する辞書領域が領域 A_1 、 A_0 、 A_2 と変移していったときの符号化データを示す。

【0055】即ち、最初は辞書領域 $j = 1$ にあることから辞書領域番号 $j = 1$ を符号化し、続いて領域 A_1 に属する文字列の符号化データを出力する。符号化データを3つ出力すると4番目の符号化データは領域 A_0 に属していることから、ここで領域 $j = 0$ を符号データとした後に文字列の符号データを出力する。更に、符号化データが領域 A_2 に属すると領域 $j = 2$ を符号化して出力した後に領域 A_2 に属する文字列の参照番号の符号データを出力する。

【0056】図12は図8のLZW符号化アルゴリズムで得られた符号データから元の文字列を復元するLZW復号化アルゴリズムを示したフローチャートである。図12において、まずステップS1で図8の符号化と同様、初期値 S_{00} 、 S_{10} 、 S_{20} をそれぞれ辞書の対応領域 A_{00} 、 A_{10} 、 A_{20} に格納する。続いてステップS2で符号を入力し、ステップS3で辞書領域番号の符号入力の有無をチェックし、辞書領域番号があればステップS5で現在の辞書領域番号を更新してステップS2で本来の符号を入力する。

【0057】続いてステップS4で現在の辞書領域に対応する番号 i_1 である符号を、前述したモード1～3の対応番号を求める関係式を使用して辞書の参照番号 i に戻す。次にステップS5で辞書を参照し、参照番号 i に対応する文字列を復元する。ステップS6で前回の辞書領域 A_k に空き領域があれば辞書アドレス n_k を1つインクリメントし、前回の復元した文字列に今回復元した文字列の先頭文字を付加した文字列を辞書領域 A_k に辞書アドレス n_k を付けて登録する。

【0058】以上の処理をステップS7で全ての符号データの入力が済むまで繰り返し、符号データの入力がなくなれば処理を終了する。図13は本発明の第2実施例で使用する辞書構成を示した説明図である。即ち、図9に示す辞書構成の実施例にあつては、辞書をデータ1、2に共通の領域 A_0 とデータ1、2に特有な領域 A_1 、 A_2 に分けていたが、図13の実施例にあつては、データ1、2に共通な領域の各々とデータ1、2に特有な領域を一緒にして1つの辞書領域としたことを特徴とす

る。

【0059】即ち、図13の辞書構成にあっては、データ1, 2に共通な高頻度の初期値 S_{00} を登録した領域については、データ1, 2に固有な領域 A_1 または A_2 の一部に含ませており、この共通の初期値 S_{00} にデータ1, 2に特有な領域 A_1 , A_2 のそれぞれを加えた領域が実際の符号化に使用するデータ1, 2に固有な辞書領域となる。

【0060】図13のように共通領域を各データに固有の領域と一緒にした場合の辞書構成におけるLZW符号化アルゴリズムは図8と同じになるが、図8のステップS6における最長一致した辞書の参照番号 i を辞書領域に対応する番号 i_1 に変換するモード1~3の処理が異なる。図13の第2実施例における参照番号 i を各領域の対応番号 i_1 に直す処理は次のモード1~3のようになる。

【モード1】

$0 \leq i < N_{0max}$; 辞書領域番号=1または2

辞書領域の対応番号 $i_0 = i + Np$

【モード2】

$N_{0max} \leq i < N_{1max}$; 辞書番号=1

辞書領域の対応番号 $i_1 = n_0 + i - N_{0max} + Np$

【モード3】

$N_{1max} \leq i$; 辞書番号=2

辞書領域に対応番号 $i_2 = n_0 + i - (N_{0max} + N_{1max}) + Np$

この第2実施例におけるモード1~3における対応番号 $i_0 \sim i_2$ への変換は、モード1の共通領域 A_0 については最初の実施例と同じであるが、モード2, 3については共通領域 S_{00} の既登録要素数 n_0 分だけ領域 A_1 , A_2 を拡張するように対応番号を求める。

【0061】図14は辞書参照番号 i が領域 A_1 に属した場合の対応番号 i_1 への変換を示したもので、想像線で示す実際の辞書番号 i に対する対応番号 i_1 を求めると、領域 A_1 の一部である共通領域 A_0 の既登録要素数 n_0 分だけ領域 A_1 を拡張した参照番号に変換することを意味する。これによって、領域 A_1 は共通領域 A_0 を含む1つの領域として扱われることになる。

【0062】その結果、入力データと最長一致する文字列が領域 A_1 または A_2 の一部である共通領域 A_0 に属するときは辞書領域番号 $j = 1, 2$ の指定は不要となる。この共通領域 A_{00} を各データ特有の領域 A_{10}, A_{20} と一緒にした辞書構成による符号化は、結局は2つの辞書を切り換えて使用していることと同じになる。また、上記の実施例におけるLZW符号化にあっては、入力する混在データの統計的性質の変動も考慮し、共通初期値 S_{00} の登録領域 A_{00}, A_{10}, A_{20} に続いて空き領域 A_{02}, A_{12}, A_{22} を設け、実際のLZW符号化で得られた新たな文字列を登録する学習領域とし、学習によって混在データの統計的性質の変動を吸収している。

【0063】しかしながら、入力データの統計的性質がデータの種類ごとに予め分かっている辞書の初期値登録領域に続く空き領域 A_{02}, A_{12}, A_{22} は設けず、初期値だけで符号化を行っても良い。このように初期値 S_{00}, S_{10}, S_{20} のみでLZW符号化を行った場合には辞書への登録操作が省略できるため、更に処理速度を向上させることができる。

【0064】更に、上記の実施例にあっては2種類のデータの符号化に適用した場合を例にとるものであったが、本発明はこれに限定されず、2種類以上のデータについても全く同様に適用することができ、この場合にはデータの種類毎に高頻度の共通部分と各データ固有の高頻度の部分とに分けて集合を作り、各集合毎に参照番号を割り振って符号化すれば良い。

【0065】

【発明の効果】以上説明したように本発明によれば、複数種類のデータについて調べた高頻度の出現文字列の初期値を1つの辞書に登録して複数種類の混在データの符号化及び復号化を行うことができ、単一辞書であることからソフトウェアによるシーケンシャル処理であってもデータの種類毎に分割辞書を用いた方法に比べ、より高速の処理を行って高圧縮率を得ることができる。

【図面の簡単な説明】

【図1】本発明の原理説明図

【図2】本発明の複数種類データのデータ圧縮方法を実施する装置構成の実施例構成図

【図3】本発明の初期値作成処理の概略を示したフローチャート

【図4】本発明における初期値作成処理の内容を示した説明図

【図5】本発明の初期値作成における符号化で得られたサンプルデータ1, 2の出現頻度を示した説明図

【図6】図3の共通部分 S_{00} の初期値作成アルゴリズムを示したフローチャート

【図7】図3の固有部分 S_{10}, S_{20} の初期値作成アルゴリズムを示したフローチャート

【図8】本発明の第1実施例におけるLZW符号化アルゴリズムを示したフローチャート

【図9】図9の符号化で使用する辞書構成の説明図

【図10】図8の領域内の番号に変換する際に使用する予約語数 Np の辞書内容を示した説明図

【図11】図8の符号化で得られる符号データの説明図

【図12】本発明の第2実施例におけるLZW符号化アルゴリズムを示したフローチャート

【図13】図12のLZW符号化で使用する辞書構成の説明図

【図14】図13における領域内の対応番号が意味する辞書領域の説明図

【図15】従来のLZW符号化アルゴリズムを示したフローチャート

【図16】従来のLZW復号化アルゴリズムを示したフローチャート

【図17】データの種類毎に調べて高頻度の文字列を分割辞書に初期登録して行うLZW符号化アルゴリズムを示したフローチャート

【図18】図17の変形を示したフローチャート

【符号の説明】

10：辞書

16：CPU

18：プログラムメモリ

20：コントロールソフト

22：符号化ソフト

24：復号化ソフト

26：出現頻度カウントテーブル

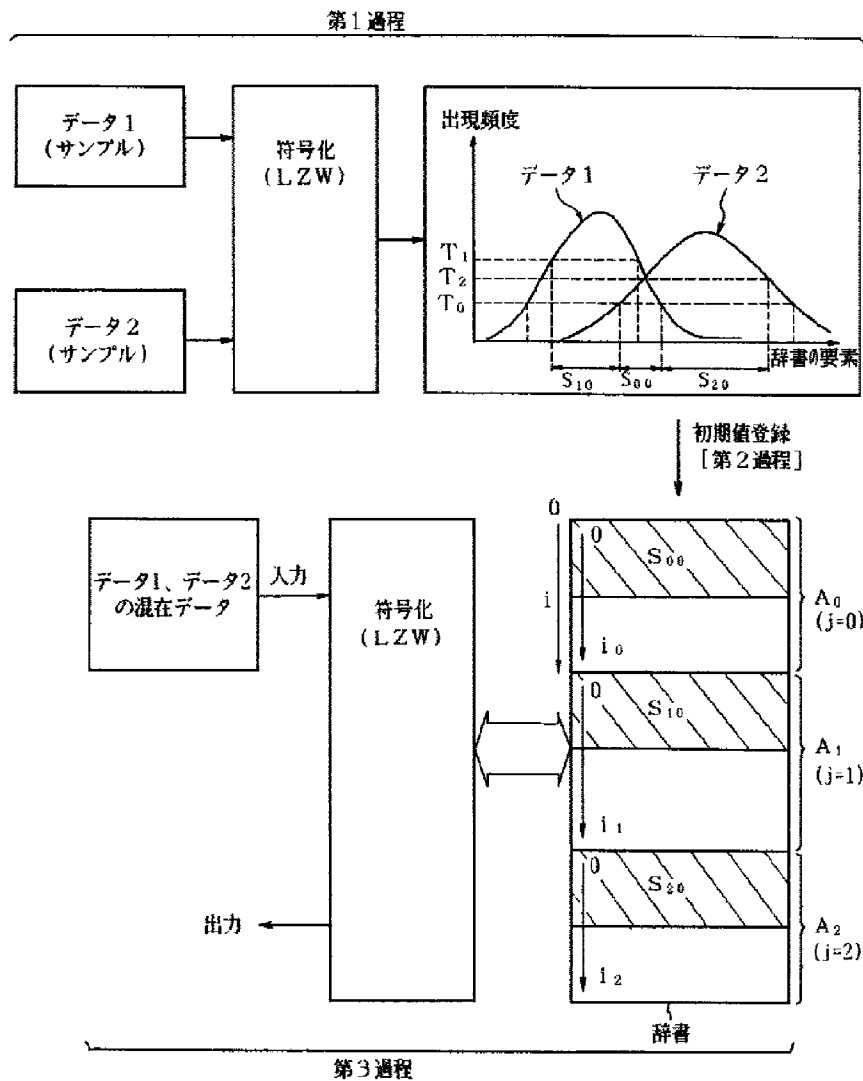
28：頻度閾値格納テーブル

30：データメモリ

32：データバッファ

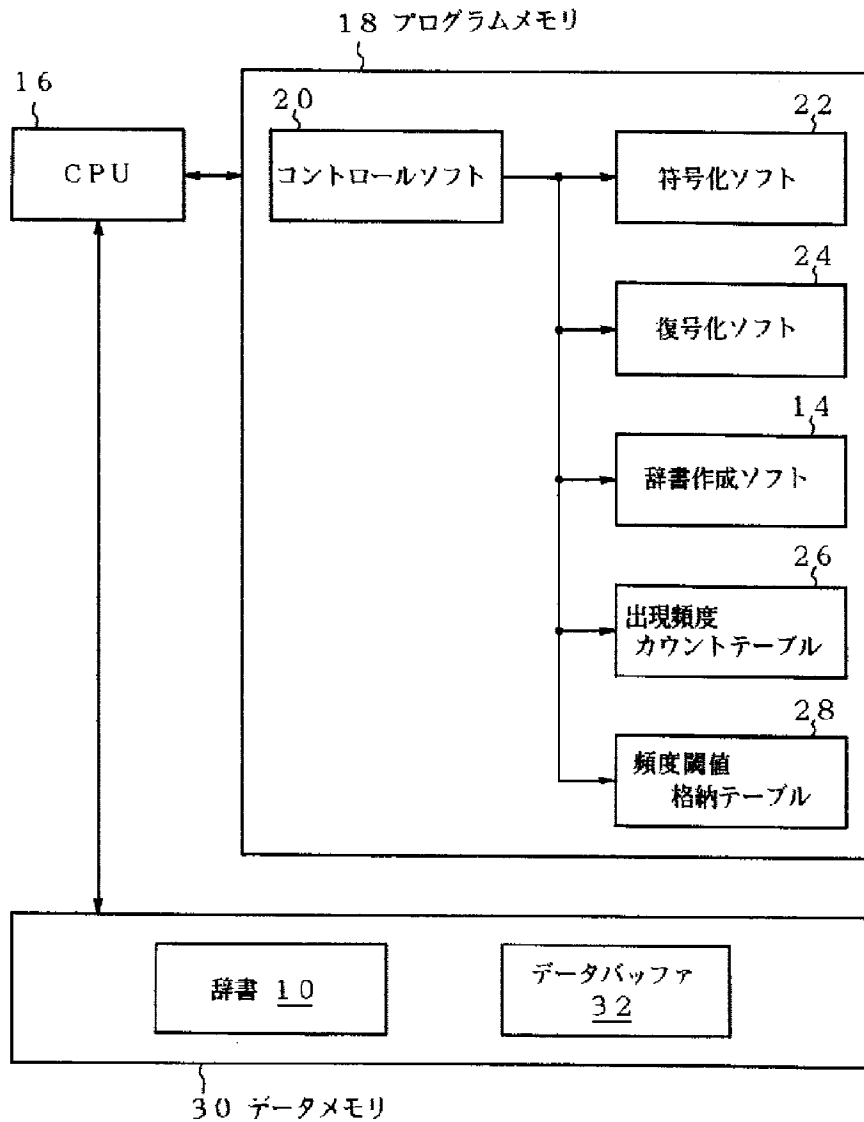
【図1】

本発明の原理説明図



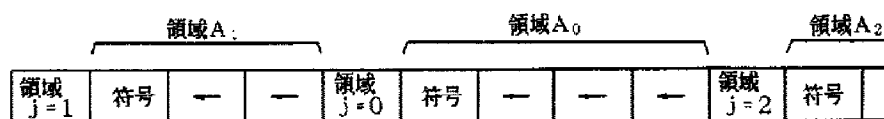
【図2】

本発明の複数種類データのデータ圧縮方法を実施する装置構成の実施例構成図



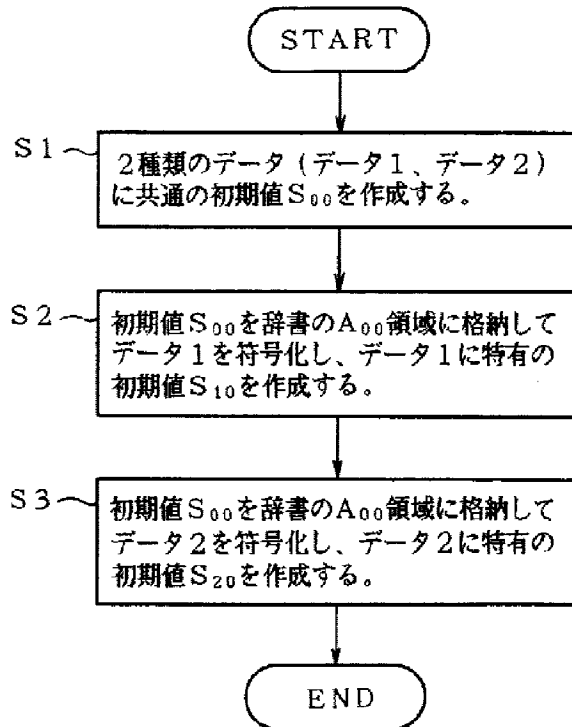
【図11】

図8の符号化で得られる符号データの説明図



【図3】

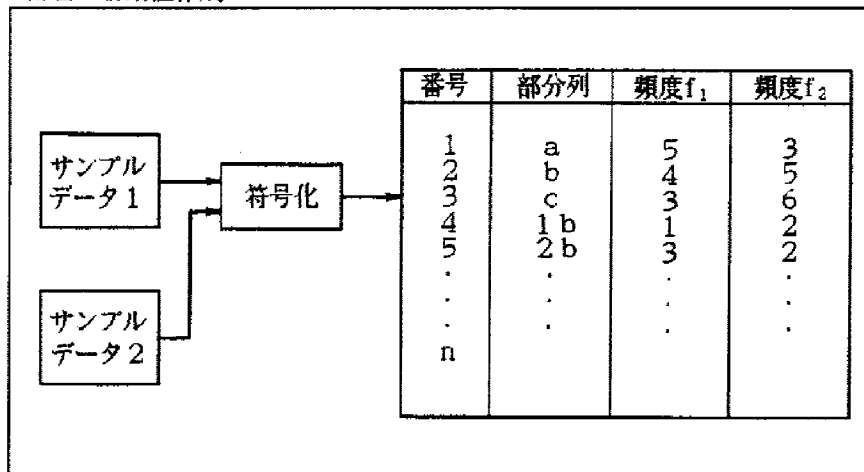
本発明の初期値作成処理の概略を示したフローチャート



【図4】

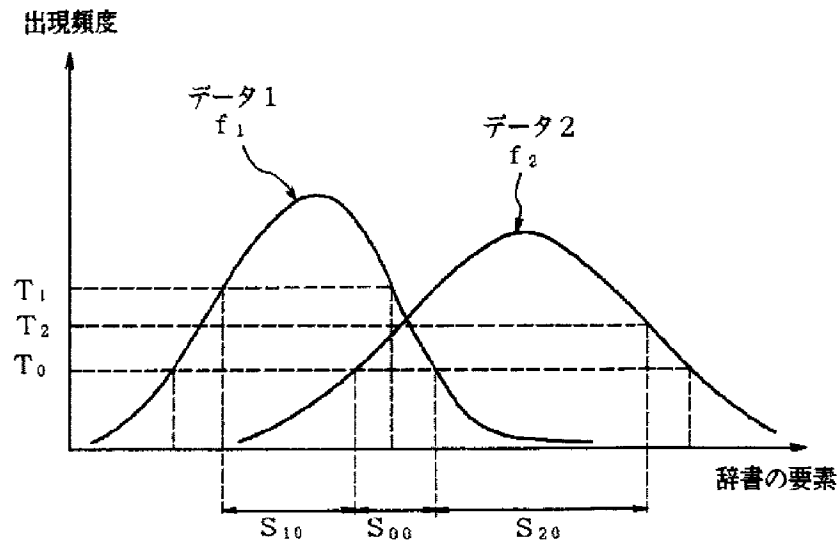
本発明における初期値作成処理の内容を示した説明図

辞書の初期値作成



【図5】

本発明の初期値作成における符号化で得られたサンプルデータ1、2の出現頻度を示した説明図

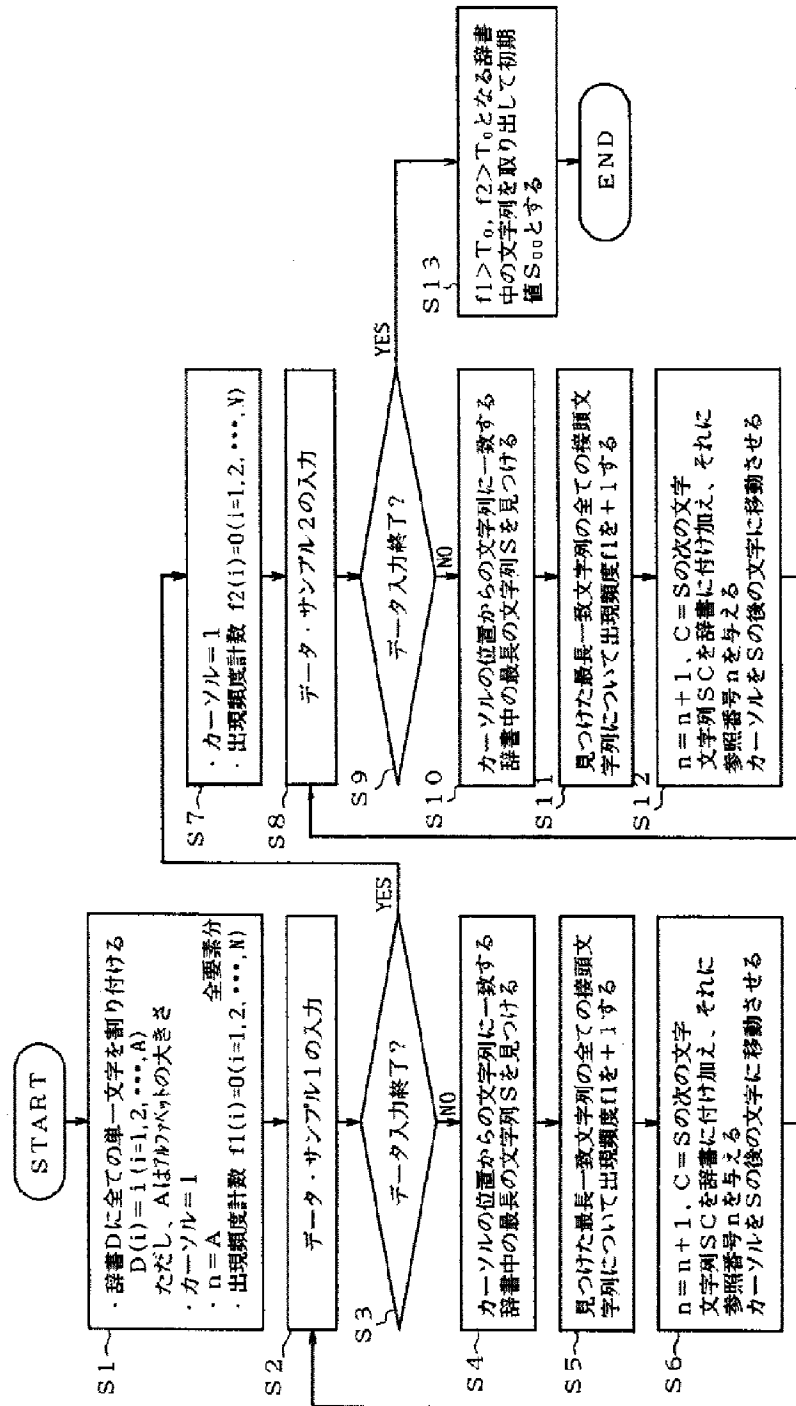


【図10】

図8の領域内の番号に変換する際に使用する予約語数Npの辞書内容を示した説明図

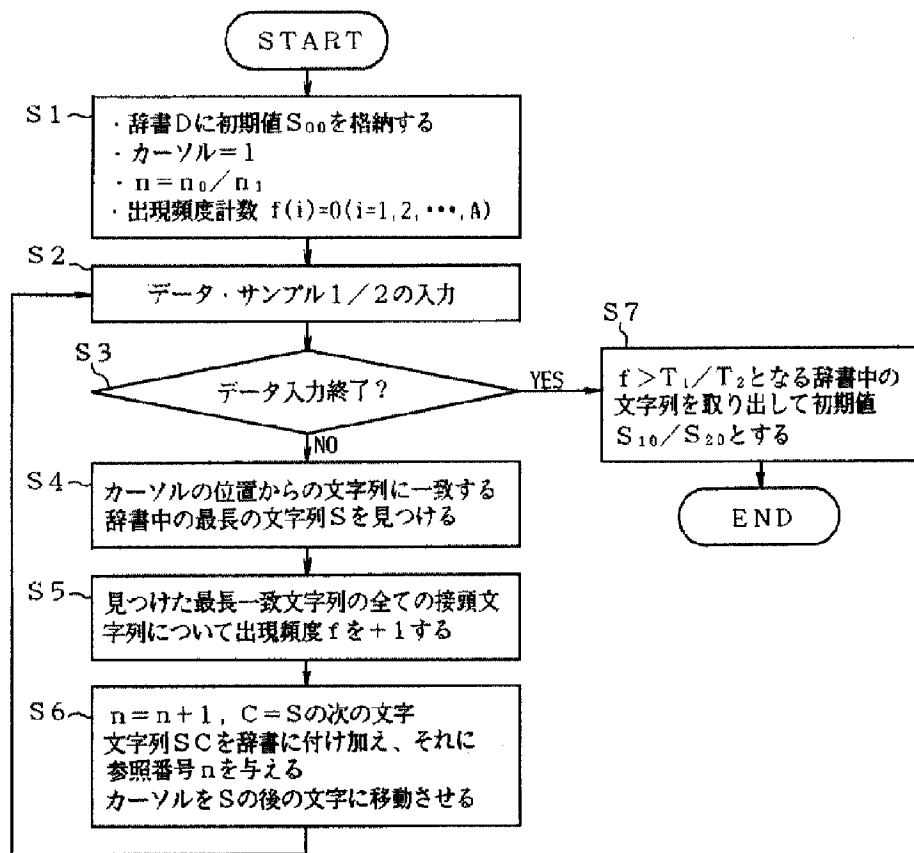
0	辞書領域番号 (A ₀)	予約語領域
1	辞書領域番号 (A ₁)	
2	辞書領域番号 (A ₂)	
3	辞書初期化コマンド	
4	EOF等	
5		辞書領域
6		
7		
8		

【図6】

図3の共通部分S₀₀の初期値作成アルゴリズムを示したフローチャート

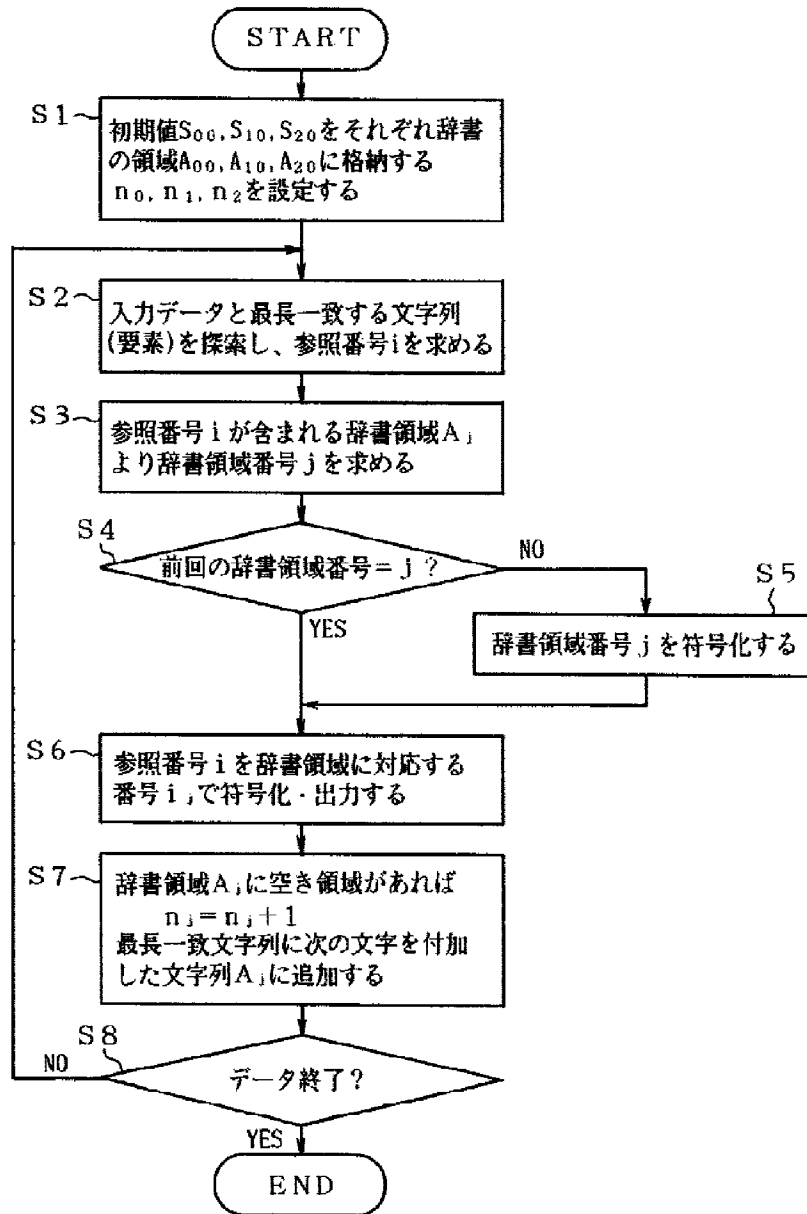
【図7】

図3の固有部分 S_{10} 、 S_{20} の初期値作成アルゴリズムを示したフローチャート



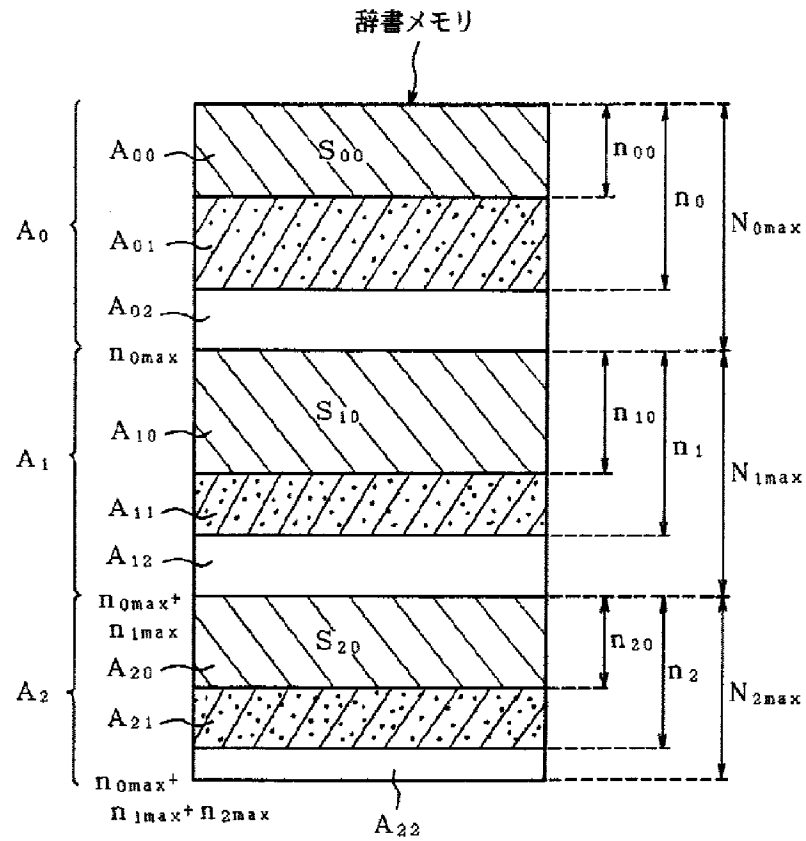
【図8】

本発明の第1実施例におけるLZW符号化アルゴリズムを示したフローチャート



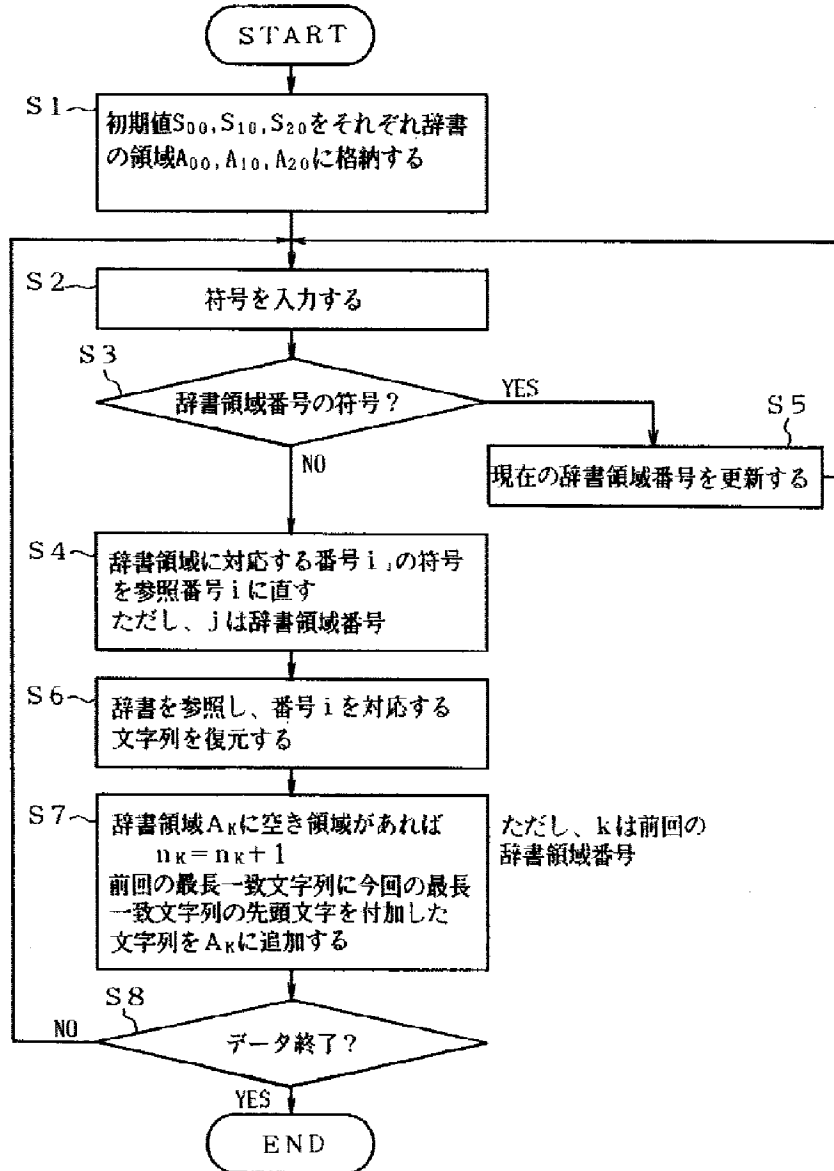
【図9】

図8の符号化で使用する辞書構成の説明図



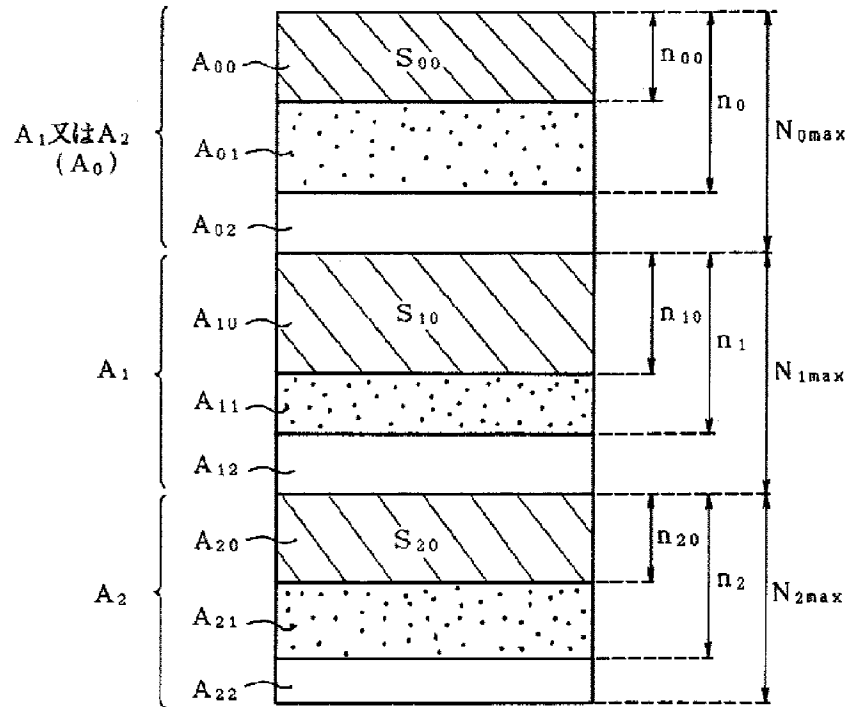
【図12】

本発明の第2実施例におけるLZW符号化アルゴリズムを示したフローチャート



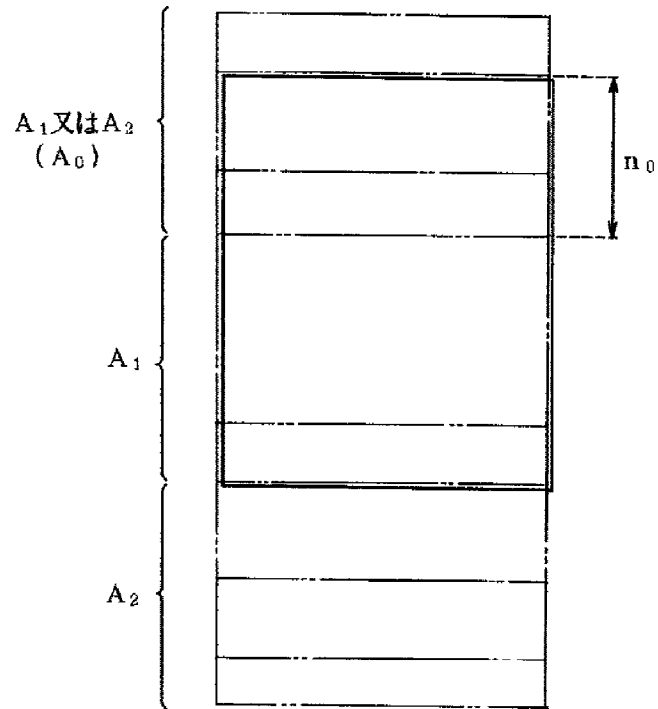
【図13】

図12のLZW符号化で使用する辞書構成の説明図



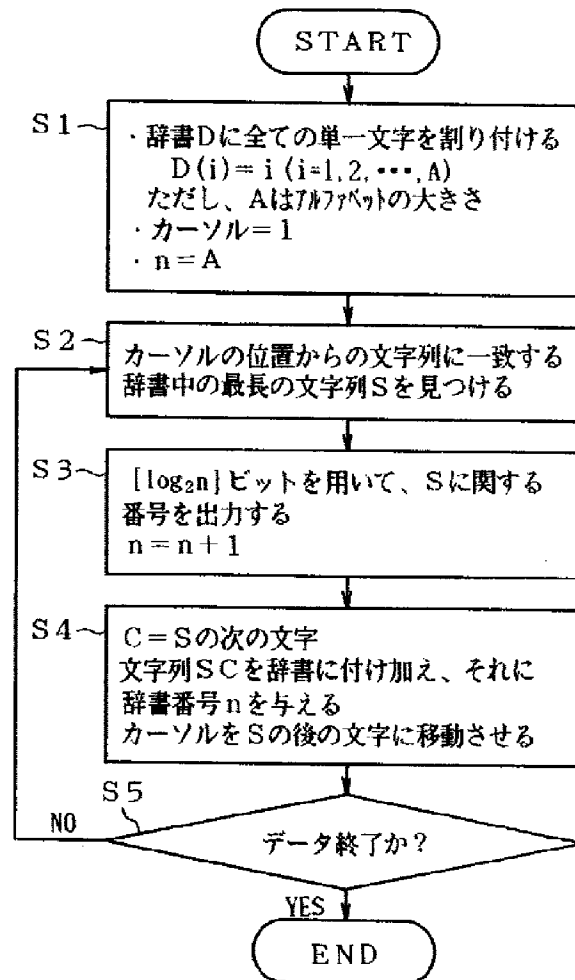
【図14】

図13における領域内の対応番号が意味する辞書領域の説明図



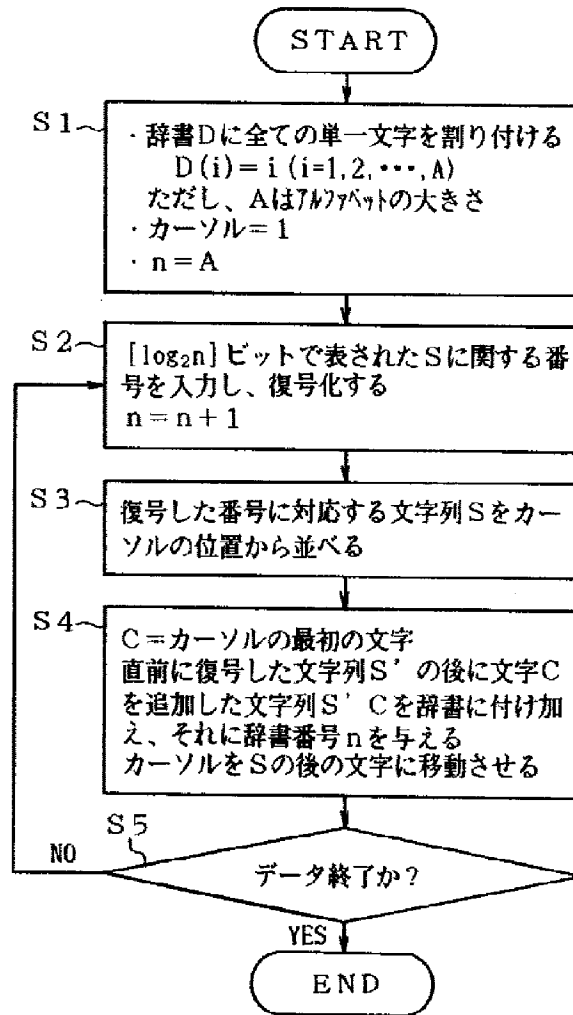
【図15】

従来のLZW符号化アルゴリズムを示したフローチャート



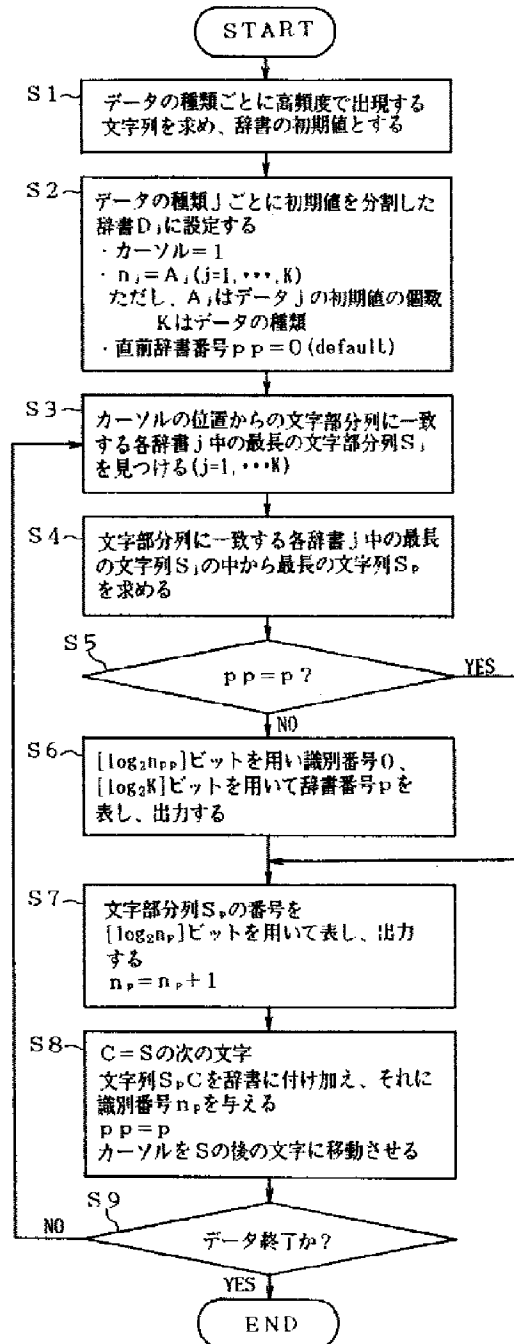
【図16】

従来のLZW復号化アルゴリズムを示したフローチャート



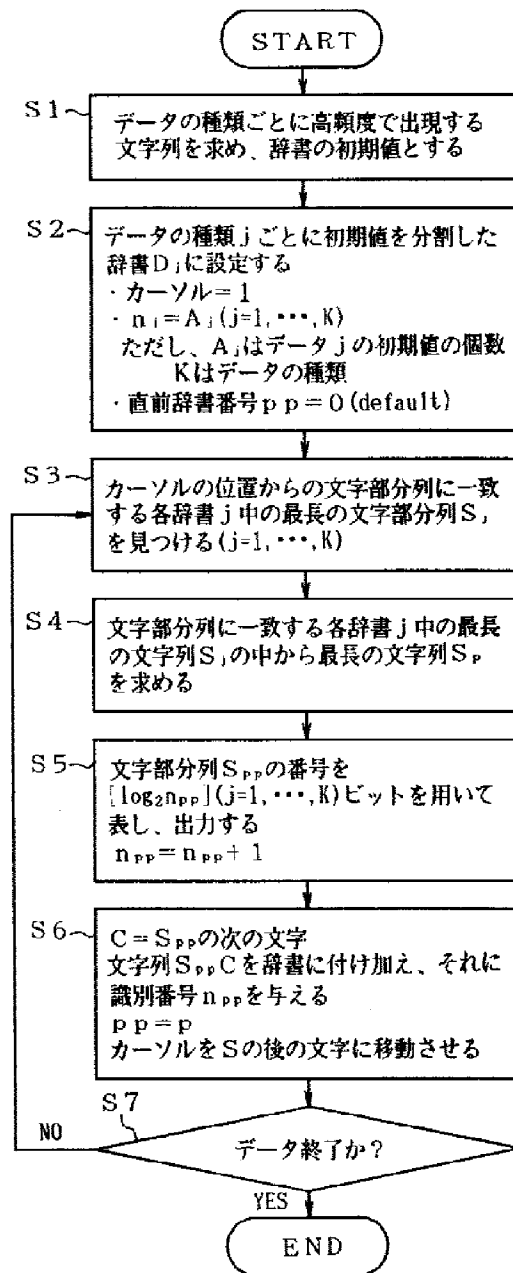
【図17】

データの種類毎に調べて高頻度の文字列を分割辞書に初期登録して行うLZW符号化アルゴリズムを示したフローチャート



【図18】

図17の変形を示したフローチャート



フロントページの続き

(72)発明者 千葉 広隆
 神奈川県川崎市中原区上小田中1015番地
 富士通株式会社内